

# SciDB I/O Performance Experiment

Parker Abercrombie  
parker@pabercrombie.com

July 22, 2014

## Experimental setup

The goal of this experiment was to determine how different storage configurations affect SciDB performance on EC2.

This experiment was performed using two m3.2xlarge EC2 VMs, each running 8 SciDB instances. One VM was EBS optimized and the other was not. SciDB data was stored on EBS volumes formatted using the ext4 filesystem. The following parameters were varied:

1. Number of attached volumes. Each volume was 200 GB.
2. Type of attached volumes. Both General Purpose SSD and Provisioned IOPS SSD were tested.
3. EBS optimized VM vs. not EBS optimized.

Queries were performed against the AVHRR GIMMS array with schema:

```
<ndvi:  int16 NULL,  
flag:   uint8>
```

```
[year   = 1981:2014, 10,    0,  
period = 1:24,      24,    0,  
y       = 0:2159,   500,   0,  
x       = 0:4319,   500,   0]
```

Benchmarks were recorded for the time required to restore the array from backup, and for several types of queries. For slice and one month aggregate, measurements were averaged over at least four runs of the query. For array restore and 16-day interval aggregation only one measurement was collected for each configuration. Measurements were collected for the following operations:

### Restore array

Restore the GIMMS array from a SciDB opaque backup file. Note: times reported are for a single run of the restore operation.

### Select one date

Select global AVHRR measurements for one 16-day interval. Result is a 2160x4320 array.

```
slice(gimms, year, <year>, period, <period>)
```

### Aggregate one month over all years

Compute average NDVI for one month at each pixel (i.e. compute average July NDVI at each pixel over all years). Result is a 2160x4320 array.

```
aggregate(gimms, avg(ndvi) as ndvi, y, x)
```

## Aggregate by 16-day interval

Compute average global NDVI for each 16-day interval in the year. Result is a 24 element vector (one number for January 1-15, one number for January 16-30, etc.)

```
aggregate(gimms, avg(ndvi), period)
```

## Results

EBS optimized instances performed slightly better than non-EBS optimized. General purpose volumes vs. Provisioned IOPS volumes show virtually no different, even if the VM is EBS optimized. Adding disks to the system reduced the array restore time slightly, but had little if any effect on query times.

Array restoration time shows considerable variability, and the lowest measurement was actually on a VM that was not EBS optimized. This operation was only conducted once for each configuration, and the variability is probably due to storage demand at the time that tests were run. However, it is surprising that the non-EBS optimized VM performed best, even if only one run.

I did not measure variance of the different runs, but times using the provisioned I/O seemed to be more consistent.

